# Appendix 1.  BLAST alignments

A.	hepsin vs spinesin:

**Blast 2 Sequences results**

PubMed    Entrez    BLAST    OMIM    Taxonomy    Structure
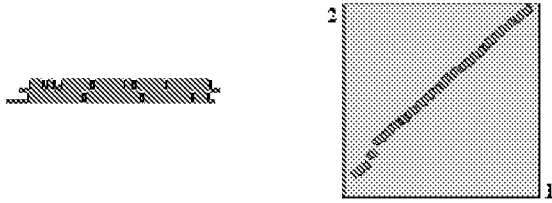
BLAST 2 SEQUENCES RESULTS VERSION BLASTP 2.2.17 [Aug-26-2007]

Matrix BLOSUM62  gap open 11  gap extension: 1
x_dropoff 0  expect 10.0000  wordsize 3  Filter  View option Standard
Masking character option X for protein, n for nucleotide  Masking color option Black
Show CDS translation

Sequence 1: gi|230057|Serine protease hepsin (Transmembrane protease, serine 1) [Contains: Serine protease hepsin non-catalytic chain; Serine protease hepsin catalytic chain]
Length = 417 (1 .. 417)

Sequence 2: gi|1234891|spinesin [Homo sapiens]
Length = 457 (1 .. 457)

NOTE.Bitscore and expect value are calculated based on the size of the nr database.

```
 Score =  221 bits (563),  Expect = 1e-55
 Identities = 140/412 (33%),  Positives = 200/412 (48%),  Gaps = 51/412 (12%)

Query  24   GTLLLLTAIGAASWAIVAVLLRSDQESLYPV-----------QVSSALARLRVFDKT---  69
                 G L LL   G SW +V  L + +S+ +A L     KT
Sbjct  53   GALGLLAGAGVGSWLLVLYLCPAASQPISGTLQDEEITLSCSEASAEEALLPALPKTVSF  112

Query  70   ---------------PGTWPLLCSSRSKARVAGLSCEEISGFLRALTHSELDVRTAGANGT  114
                            + N L+C   + +    C +G LR   H +++     N +
Sbjct  113  RINSEDFLLEAQVRDQPRNLLVCHEGWSPALGLQICWSLGHLRLTHHKFVNLTDIKLNSS  172

Query  115  SGPFCVDEGRLPHTQRLLEVI--EVCDCPRGRFLAAICQDCGRRKLFVDPIVGGRDTSLG  172
                 F +     P   LE       +C  G+ ++  C +CG R L    PIVGG+ + G
Sbjct  173  QEFRQLS----PRLGGFLEEAMQPRRNCTSGQVVSLRCSECGARPL-ASRIVGGQSVAFG  227

Query  173  RWPWQVSLRYDGAHLCGGSLLSGDWVLTAAHCFBE-RNRVLSRWRVFAGAVAQAS--SHG  229
                 RWPWQ S+    H CGGS+L+  WV+TAAHC   R   LS WRV AG V+ ++   PH
```
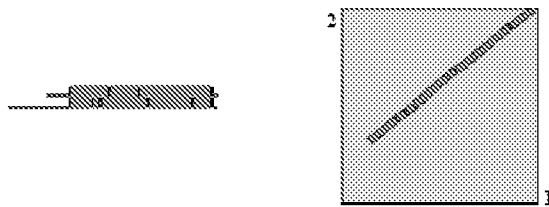
```
Sbjct  228  RWFWQASVALGFRHTCGGSVLAPRNVVIAAHTKHSFRLABLSSNRVHAGLVSHSAVRSHQ  287

Query  230  LQLGVQAVVYHGGYLSFRDSNSEENSNDIALVHLSSPLFLTEYIQPVCLPAAGQRLVDGK  289
              L V+ ++ H  Y       +++ +  D+AL+ L + L ++ + VCLPA  Q     G
Sbjct  288  GAL-VERIIPHPLY------SAQNHDYDVALLRLQTALNFSDTVGAVCLPAKEQKFPKFS  345

Query  290  ICTVTGWGNTQ-YYGQQAGVLQBARVPIISNDVCNGADFYGNQIKPKMFCAGYPEGGIDA  349
              C V+GWG+T  +  + +LQ+  VP+ S +CN +  Y   + P+N CAGY +G  DA
Sbjct  341  RCNVSGWGHTHPSHTYSSIMLQETVVPLFSTQLCNSSCVYSGALTPRMLCAGYLDGPAIA  400

Query  349  CQGDSGGSFVCEDSISRIPRNRLCGIVSWGTGCALAQKPGVYTKVSDFRENI  400
             CQGDSGGP VC D        NRL G+VSWG  CA     PGVY KV++F +NI
Sbjct  401  CQGDSGGPLVCPDG----DTWRLVGVVEWDRACAEPUHPGVYAKVAEFLDWI  449
```

16

## B.    hepsin vs TMPRSS2:

**Sequence 1**: gi|123037|Serine protease hepsin (Transmembrane protease, serine 1) [Contains: Serine protease hepsin non-catalytic chain; Serine protease hepsin catalytic chain]
Length = 417 (1 .. 417)

**Sequence 2**: gi|14602419|transmembrane protease, serine 2 [Homo sapiens]
>gi|5001734|gb|AAD37117.1|AF123453_1 transmembrane serine protease 2 [Homo sapiens]
Length = 492 (1 .. 492)



NOTE:Bitscore and expect value are calculated based on the size of the nr database.



```
 Score =  210 bits (534),  Expect = 3e-62
 Identities = 117/346 (33%), Positives = 174/346 (50%), Gaps = 17/346 (4%)

Query   63   LNVFDKTEGTWPLLCSSRSKRVAGLSCEEMGFLRALTHSELDVRTAGANGTSGFFCVDE   122
             L ++     +W  +C      N        +C +MG+      S+   V    G+
Sbjct   158  LQMYSSQRKSWHPVCQDDWNENYGRAACRDMGYKNNFYSSQGIVDDSGSTSFMKLN-TSA  216

Query   123  GRLPHTQRLLEVISVCDCPRGRFLAAICQDCG--RRKLFVDRIVGGRDTSLGRWPWQVSL  180
             G +  ++L    +   C      ++  C  C       RIVGG      G WPWQVSL
Sbjct   217  GSVDIYNELYHSDA---CSSKAVVSLRCIACDYNLNSSRQSRIVGGESALPGANPWQVSL  273

Query   181  RYDGAHLCGGSLLSGDWVLTAAHCFPERNRVLSRWRVFAGAVAQASP-HGLQLGVQRVVY  239
             +++GG+++ +S++ +AAHC +          W  FAG + Q+   +G    V+ V+
Sbjct   274  HVQNVHVCGGSIITPEWIVTAAHCVEKPLKNPWHVTAFGILRQSFNFYGAGYQVEKVIS   333

Query   240  HGGYLPFRDPNSEENSNDIALVHLSSPLPLTEYIQPVCLPAAGQALVDGKICTVTGWGNT  299
             H Y      +S+  +NDIAL+ L  PL   + ++PVCLP G L   ++C ++GWG T
Sbjct   334  RPNY------DSKTKNSDIALMFLQKPLTFNDLVKPVCLPNPNRMLQPEQLCWISGWGAT  387

Query   300  QYYGQQAGVLQEARVPIISNDVCNGADFYGNQIKPRMFCAGYPEGGIDACQGDSGGPFVC  359
             + G+ + VL +A+V +I   CN   Y N+I P N CAG+ +G +D+CQGDSG  V
Sbjct   395  EENGVTSEVLNAAKVLLIETQRCNSRYVYDNLITPAMICAGFLQGNVDSCQGDSGGPLV-  446

Query   360  EDSISRTPRWRLCGIVSWGTGCALAQKPGVYTKVSDFREWIFQAIK  405
             +S+  W L G  SWG+GCA A +PGVY   V  F +WI++  +
Sbjct   447  ---TSKNNIWWLIGDTSWGSGCAKAYRPGVYGNVMVFTDWIYRQMR  489
```
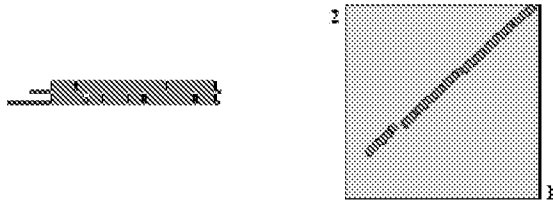
## C.  hepsin vs TMPRSS3:

Sequence 1: gi|123057|Serine protease hepsin (Transmembrane protease, serine 1) [Contains: Serine protease hepsin non-catalytic chain; Serine protease hepsin catalytic chain]
Length = 417 (1 .. 417)

Sequence 2: gi|37182183|TMPRSS3 [Homo sapiens]
Length = 432 (1 .. 432)

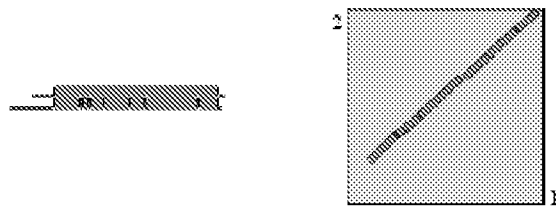NOTE.Bitscore and expect value are calculated based on the size of the nr database.

```
 Score =  221 bits (564),  Expect = 9e-56
 Identities = 125/360 (35%), Positives = 175/360 (48%), Gaps = 34/360 (9%)

Query   50   PLYPWQVSSADARLGVFDYTEGTWRLLCSSRSNARVAGLSTEENGFLPALT---HSELDV   106
             P  V++S   + L V D   G N   C      +A  +C +NG+ PA+      +LDV
Sbjct  105   PAVAVRLSNDRSTLQVLNSATGNNFSACFDNFTEALAETACPQMGYSRAVEIGPDQDLDV   155

Query  107   RTAGANGTSGFFCVDEGRLFEDQRLLEVISVCDCPRSRFLAAICQDCGRRNLPVDRIVGG   166
                       N           +Q L    S   C G ++  C CG+  L    R+VGG
Sbjct  160   VEITEN---------------EQELPMPNSSGPCLSGSLVSLHCLACGN-SLNTPRIVGG   203

Query  167   SDTSLDPNPWQVSLRYDGAHLCGGSLLSGDNVLTAAHCFPEPNRVLSRNRVFAGAVAQAS   226
             + S+ NPWQVS++YD H+CGGS+L   NVLTAAHCF + V + W+V AG+   S
Sbjct  204   EEASVDSNPWQVSIQYDKQHVCGGSILDPHWVLTAAHCFRKHTDVFN-WKVRAGSDKLGS   262

Query  227   PHSLQLGVQAVVYHGGYLPFRDPNSEENSNDIALVHLSSPLPLTEYIQPVCLPAAGQALV   286
              L +    ++       P            NDIAL+ L  PL +  ++P+CLP  + L
Sbjct  263   PPSLAVAKIIIIEFNPMYP--------KDMDIALNKLQFPLTFSGTVRPICLPFFDEELT   314

Query  287   DGKICTVDGWGNT-QYTGQQAGVLQEARVPIISKEVCNGADFYGNQINEKMFCAGYPEGG   345
              + GWG T Q  G+ + +L +R V +I +  CN  D Y  ++  KM CAG PEGG
Sbjct  315   PATPLNIIGWGFTKQNGGKMSDILLQAEVQVIDSTRCNADDAYQGEVIEKNGCAGIPEGG   374

Query  346   IDACQGDSGGPFVCEDSISRTPRRRLQGIVSWGTGCALAQKPGVYTKVSDFREWIFQAIK   405
             +D CQGDSGGP + +     + +M + GIVSWG GC      PGVYT+VS + WI+    N
Sbjct  375   VDTCQGDSGGPLMTQ-----SDQNEVVGIVSWGYGCGGPEIPGVYTRVSAYLNWIYRVRK   429
```

## D.  hepsin vs TMPRSS4:

NOTE:Bitscore and expect value are calculated based on the size of the nr database.

```
 Score =  215 bits (558),  Expect = 4e-55
 Identities = 126/367 (34%), Positives = 180/367 (50%), Gaps = 23/367 (6%)

Query   50   PLYPWQVSSADASLMVFDKTEGTWRLLCSSRSNARVAGLSCEEMGFLRALTHSELDVKTA  109
             P V++S  + L V D  G W C      +A +C +MG+    T    ++
Sbjct  100   PAVAVSLSMDRSILQVLDSATGNWFSACFDNFTEALAETACRQMGYSSKPTFRAVEI---  156

Query  115   GANGTSGFFCVDEGRLPHIQRLLEVISVCDCPRGRFLAAICQDCGPRELPVDRIVGGRDT  169
             G +     + E  ++Q L   S  C  G  +    C  CG+  L    R+VGG +
Sbjct  157   GPDQDLDVVEITE----NSQELPMRNSSGPCLSGSLVSLHCLACGK-SLNTPRVVGGEEA  211

Query  170   SLGSWPWQVSLRYDGAHLCGGSLLSGKWVLTAAHCFPERNRVLSRWRVFAGAVAQASPHG  229
             S+  WPWQVS++YD  H+CGGS+L   WVLTAAHCF +  V + W+V AG+    S
Sbjct  212   SVDSWPWQVSIQYDNQHVCGGSILDPHWVLTAAHCFRKHTIDVFN-NIVRAGSDNLGSFPS  270

Query  230   LQLGVQANVYHGGYLPFRDPNSEENSNDIALVHLSSPLPLTEYIQPVCLPAAGQRLVDGK  289
             L +   ++     P        NDIAL+ L  PL +  ++P+CLP   + L
Sbjct  271   LAVAKIIIIEFNPMYP--------SDMDTALMPLQFPLTFSGTVRPICLPFFDEELIPAT  322

Query  290   ICTVTGWGNT-QIYGQQAGVLQEARVPIISNDVCNGADFYGNQIKPKMFCAGYPEGGIDA  348
              + GWG T Q  G+ + +L +A V +  CN  D Y  ++ PM CAG PEGG+D
Sbjct  323   PLSIIGWGFTKQNGGKMSDILLQASVQVIDSTRCNADDAYQGEVTEKMMCAGIPEGGVDT  382

Query  349   CQGDSGGPFVCEDSISRTPRWRLCGIVSWGTGCALAQKPGVYTKVSDFREWIFQAIK  405
             CQGDS GP + +      + +N + GIVSWG GC     PGVYTKVS  +I+  K
Sbjct  383   CQGDSGGPLNIQ-----SDQNVFGIVSWGYGCGGPSTPGVYTEVSAYLNWIINTWK  434
```

## E.    hepsin vs enteropeptidase:

Sequence 1. gi|11307|Serine protease hepsin (Transmembrane protease, serine 1) [Contains: Serine protease hepsin non-catalytic chain; Serine protease hepsin catalytic chain]
Length = 417 (1 .. 417)

Sequence 2. gi|6650008|enteropeptidase [Homo sapiens]
Length = 1019 (1 .. 1019)



NOTE Bitscore and expect value are calculated based on the size of the nr database.



```
 Score =  211 bits (536),  Expect = 1e-52
 Identities = 108/256 (40%), Positives = 146/256 (56%), Gaps = 17/256 (6%)

Query  151   QDCGRRKLPVD---RIVGGRDTSLGRWPWQVSLRIDGAHLCGGSLLSGRWVLTAAHCFPE   207
             +  CG++       D    +IVGG +    G NPW V L Y G   LCG SL+S  N+++AAHC
Sbjct  770   KSCGKKLAAQDITPDIVGGSNAKEGAWPWVVGLYYGGRLLCGASLVSSIWLVSAAHCVYG   828

Query  205   PHPNLSPNNVFAGAVAQA---EPHGLQLGVQATVTHGGYLSFRDPNEENSNDIALVHLS   264
             R          G  ++   EP +  + +V + Y       N    NDIA++KL
Sbjct  830   RNLEPSKWIAILGLHMKSNLTSPQTVPRLIDEIVINPHY------NPPRKDNDIANMHLE   883

Query  265   SPLPLTEYIQPVCLPAAGQALVDGKICTVTGWGNTQYYGQQAGVLQEARVPIISNDVCNG   324
               +  T+YIQP+CLP   Q   G+ C++ GWG   Y G  A +LQEA VP++SN+ C
Sbjct  884   FKVDYTDYIQPICLPEENQVFPPGRSCSIAGWGTVVYQGTTANILQEADVPLLSNERCQ-   942

Query  325   ADFYGNQIKPKMFCAGYPEGGIDACQGDSGGPFVCEDSISRTPRWRLCGIVSWGTGCALA   384
             I    N CAGY EGGID+CQGDSGGP +C+++     RW L G+ S+G  CAL
Sbjct  943   QQMPEYNITENMICAGYEEGGIDSCQGDSGGPLMDQEN----NRWFLAGVTSFGYKCALP   998

Query  385   QKPGVYTKVSDFREWI   400
             +PGVY +VS F INI
Sbjct  999   NRPGVYARVSRFTEWI   1014
```
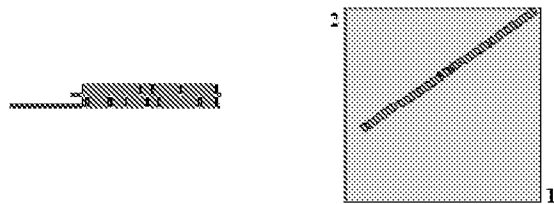
## F.  hepsin vs MSPL:

Sequence 1: gi|133057|Serine protease hepsin (Transmembrane protease, serine 1) [Contains: Serine protease hepsin non-catalytic chain; Serine protease hepsin catalytic chain]
Length = 417 (1 .. 417)

Sequence 2: gi|116256363|transmembrane protease, serine 13 [Homo sapiens]
Length = 567 (1 .. 567)



NOTE:Bitscore and expect value are calculated based on the size of the nr database.



```
 Score =  231 bits (589),  Expect = 1e-59
 Identities = 134/372 (36%),  Positives = 196/372 (52%),  Gaps = 28/372 (7%)

Query  38   IVAVLLRSEQEPLIPVQVSSADARLMVFEKTEGTWRLLCSSRSNARVAGLSCEEMGFLRA   98
            +V  L+SE+  L  V+      + L ++  +   W  +CSS N   +  +C+++GF  A
Sbjct  214  VVDINLKSDE--LGCVRFINDKSLLKIYSGSSHQWLPICSSNWNDSYSEKPCQQLGFESA   271

Query  99   LTHSELEVRTAGANGTSGFFCVDEGRLPSTQRLLEVISVCDCPRGRFLAAICQDCGRRNL  169
            + E+  R        F     L +   E +   +CP  +++   C  CG R +
Sbjct  272  KRTTEVAHRD---------FANSFSILRYNSTIQEELHREECPSQRYISLQCSHCGLRAN  322

Query  159  PVERIVGGRDTSLGRWPWQVSLRIDGAHLCGGSLLSGNWVLTAAHC-FPERNRVLSRWRV  217
               RIVGG  S +WPWQVSL +  H+CGG+L+   WVLTAAHC P R +VL  W+V
Sbjct  323  -TGPIVGGALASDSKWPWQVSLHFGTTHICGGTLIDAQWVLTAAHCFFVTREKVLEGWKV  381

Query  218  FAGAVAQASPHGL--QLGVQANVYHGGYLPFRDPNSEENSNDIALVHLSSPLPLTEYIQB  275
            +AG    ++ H L    + ++ + V      EE+  DIAL+ LS PL L+ +I P
Sbjct  382  YAGT---SNLEQLPEAASIAEIIINSNY------TDEEDDYDIALMRLSNPLTLSAHIHP  432

Query  276  VCLPAAGQALVDGKICTVRSWGNPQYYGQQAG-VLQEARVPIISNEVCNGADFYGNQIKP  334
             CLP  GQ    +  C +R+WG+ +      + L+E +V +I   CN   Y  + P
Sbjct  433  ACLPMHGQTFSLNEICWITPEPNTRETDDKTSPFLREVQVRLIDEHYCNDYLVYDSYLTP  491

Query  336  KMFCAGYPEGGIDACQGDSGGPFVCEDSISRTPRWRLCGIVSWGTGCALAQKPGVYTKVS  394
            +M CAG   GG D+CQGDSGGP VCE    + PR L G+ SWGTG   KPGVYTKV+
Sbjct  492  SMMCAGDLRGGRDSCQGDSGGPLVCE----QNNRWYLAGVTSWGTGCSQRNHPGVYTKVT  548

Query  395  DFREWIFQAIKT  406
            +   WI+  +++
Sbjct  549  EVLPWIYSMMEG  560
```